

Graph sequences with finite-time convergence for decentralized average consensus and applications in distributed optimization

Xin Jiang

School of Operations Research and Information Engineering
Cornell University

Cornell Young Researchers Workshop

October 10, 2024

Decentralized average consensus

Classic problem setup

- a **connected** graph/network $G = (V, E)$ with n agents
- each agent initially holds a vector $x_i \in \mathbb{R}^d$
- each agent only communicates with its neighbors (message passing)
- a round of communication is represented as **matrix–vector product**

$$X^{(k+1)} = WX^{(k)}, \quad \text{where } X^{(0)} = [x_1 \ x_2 \ \cdots \ x_n]^T \in \mathbb{R}^{n \times d}$$

where $W \in \mathbb{S}^n$ is the **mixing matrix**: $W_{ij} = 0$ if $\{i, j\} \notin E$

Goal: via rounds of communication, without a central agent, all agents obtain

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Classic result: **asymptotic** convergence for any $\{x_i\}$ if and only if

$$W\mathbf{1} = \mathbf{1}, \quad W^T\mathbf{1} = \mathbf{1}, \quad 1 = |\lambda_1| > |\lambda_2| \geq \cdots \geq |\lambda_n|$$

Modern applications in distributed optimization

in traditional applications and federated learning

- agents are connected with low-bandwidth channels
- communication is highly fragile; occasional link failures
- network topology is fixed or cannot be controlled

we consider modern scenarios with high-performance [data-center](#) clusters

- all GPUs are connected with high-bandwidth channels
- communication is highly reliable; no occasional link failure
- network topology can be fully controlled

Graph sequence with finite-time consensus property

the **finite-time consensus** property is defined for a given sequence of graphs

$$\{G^{(l)} \equiv (V, W^{(l)}, E^{(l)})\}_{l=0}^{\tau-1}$$

sparsity is desirable: each graph $G^{(l)}$ might not be **connected**

Consensus perspective: decentralized averaging converges in τ iterations

$$X^{(\tau)} = W^{(\tau-1)}W^{(\tau-2)} \dots W^{(1)}W^{(0)}X^{(0)} = \mathbb{1}\bar{x}^T$$

Matrix perspective:

$$W^{(\tau-1)}W^{(\tau-2)} \dots W^{(1)}W^{(0)} = \frac{1}{n}\mathbb{1}\mathbb{1}^T =: J$$

Preview

we study three classes of graph sequences with finite-time consensus

graph sequence	size n	τ
p -peer hyper-cuboids [NJYU'23]	any $n \in \mathbb{N}_{\geq 2}$	# prime factors
SDS factor graphs [JNUY'24]	any $n \in \mathbb{N}_{\geq 2}$	flexible*
DSHB factor graphs [JNUY'24]	any $n \in \mathbb{N}_{\geq 2}$	flexible*

and their applications in distributed optimization algorithms

SDS: sequential doubly stochastic; DSHB: doubly stochastic hierarchically banded

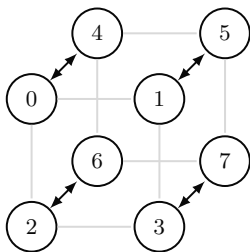
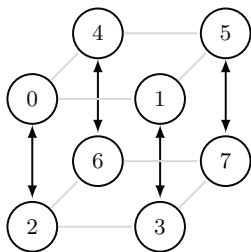
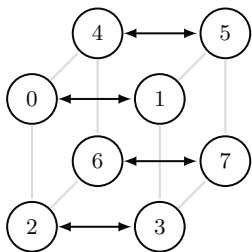
*: τ is related to a partition $n = \sum_{k=1}^{\tau} n_k$

References

- [NJYU'23] On graphs with finite-time consensus and their use in gradient tracking, *arXiv:2311.01317*; under 2nd round review in SIOPT
- [JNUY'24] Sparse factorization of the square all-ones matrix of arbitrary order, *arXiv:2401.14596*; under 2nd round review in SIMAX

One-peer hyper-cube (for $n = 2^T$)

binary representation of integers: $n = 8 = 2 \times 2 \times 2$



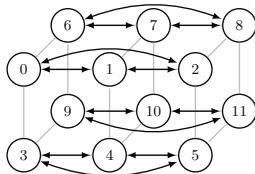
$$\begin{bmatrix} \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 & 0 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 & 0 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 & 0 \\ 0 & 0 & \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 & 0 \\ 0 & 0 & \frac{1}{2} & \frac{1}{2} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \frac{1}{2} & \frac{1}{2} & 0 & 0 \\ 0 & 0 & 0 & 0 & \frac{1}{2} & \frac{1}{2} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

$$\begin{bmatrix} \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 & 0 & 0 & 0 \\ 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 & 0 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 & 0 & 0 & 0 \\ 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ 0 & 0 & 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ 0 & 0 & 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ 0 & 0 & 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ 0 & 0 & 0 & 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & 0 & 0 & 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \end{bmatrix}$$

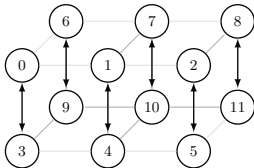
$$\begin{bmatrix} \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} & 0 & 0 & 0 \\ 0 & \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} & 0 & 0 \\ 0 & 0 & \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} & 0 \\ 0 & 0 & \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} & 0 \\ \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} & 0 & 0 & 0 \\ 0 & \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} & 0 & 0 \\ 0 & 0 & \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} & 0 \\ 0 & 0 & \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} & 0 \\ 0 & 0 & 0 & \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} \\ 0 & 0 & 0 & \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} \end{bmatrix}$$

p -Peer hyper-cuboid: An example $n = 12$

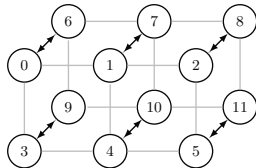
multi-base representation of integers: $n = 12 = 3 \times 2 \times 2$



$G^{(0)}$



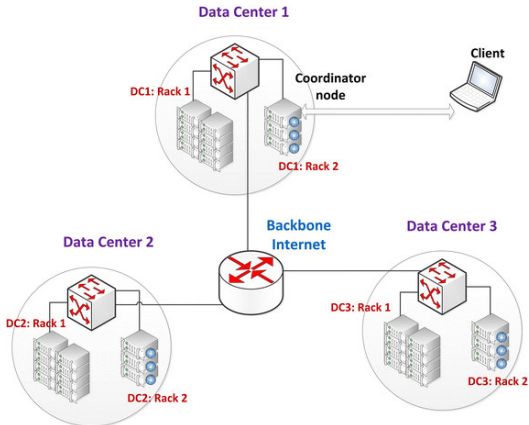
$G^{(1)}$



$G^{(2)}$

p -Peer hyper-cuboid: Limitations

- p -peer hyper-cuboids revert to fully-connected graphs when n is *prime*
- data centers are not equidistant but formed in clusters
 - intra-cluster communication is cheap, flexible and can be varied
 - inter-cluster communication is expensive and should be minimized



Three-phase communication protocol

- phase 1: intra-cluster communication achieving finite-time consensus
- phase 2: **limited** inter-cluster communication
- phase 3: intra-cluster communication achieving finite-time consensus

we now focus on reducing the communication cost in phase 2

A two-block example

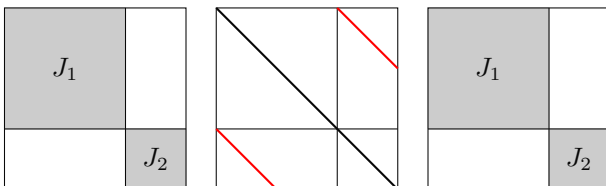
$$J = \begin{bmatrix} J_1 & \\ & J_2 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} \\ A_{12}^T & A_{22} \end{bmatrix} \begin{bmatrix} J_1 & \\ & J_2 \end{bmatrix} = \begin{bmatrix} J_1 A_{11} J_1 & J_1 A_{12} J_2 \\ (J_1 A_{12} J_2)^T & J_2 A_{22} J_2 \end{bmatrix},$$

where $n = n_1 + n_2$ with $n_1 \geq n_2$, $J_1 = \frac{1}{n_1} \mathbf{1}_{n_1} \mathbf{1}_{n_1}^T$, and $J_2 = \frac{1}{n_2} \mathbf{1}_{n_2} \mathbf{1}_{n_2}^T$

J_1		A_{11}	A_{12}	J_1	
	J_2	A_{12}^T	A_{22}		J_2

A two-block example

$$J = \begin{bmatrix} J_1 & \\ & J_2 \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} \\ A_{12}^T & A_{22} \end{bmatrix} \begin{bmatrix} J_1 & \\ & J_2 \end{bmatrix} = \begin{bmatrix} J_1 A_{11} J_1 & J_1 A_{12} J_2 \\ (J_1 A_{12} J_2)^T & J_2 A_{22} J_2 \end{bmatrix}$$



$$A = \left[\begin{array}{cc|c} \frac{n_2}{n} I_{n_2} & 0 & \frac{n_1}{n} I_{n_2} \\ 0 & I_{n_1-n_2} & 0 \\ \hline \frac{n_1}{n} I_{n_2} & 0 & \frac{n_2}{n} I_{n_1} \end{array} \right]$$

The general case

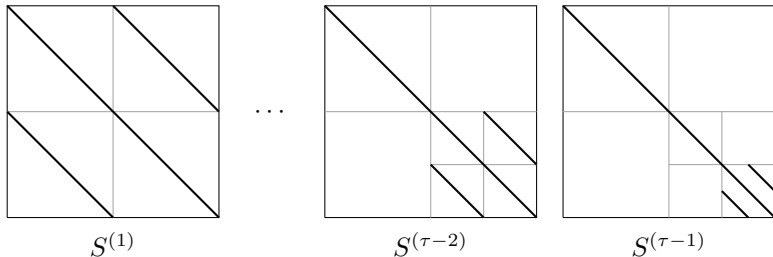
$$J = J_0 A J_0$$

- $J_0 := J_1 \oplus \cdots \oplus J_\tau$ is block diagonal with $J_k := \frac{1}{n_k} \mathbf{1}\mathbf{1}^T \in \mathbb{R}^{n_k \times n_k}$
- \oplus the direct sum of two matrices: $X \oplus Y = \text{blkdiag}(X, Y)$
- each J_k can be further decomposed into, e.g., p -peer hyper-cuboids
- key **trade-off**:
communication per round (sparsity) v.s. # communication rounds
- we provide two options for the A -factor
 - A can be decomposed as product of several banded matrices
 - A can be *hierarchically* partitioned as banded matrices

Sequential doubly stochastic (SDS) factorization

$$\begin{aligned} J &= J_0 A_L J_0 && \text{with } A_L = S^{(1)} S^{(2)} \dots S^{(\tau-1)} \\ J &= J_0 A_R J_0 && \text{with } A_R = S^{(\tau-1)} S^{(\tau-2)} \dots S^{(1)} \end{aligned}$$

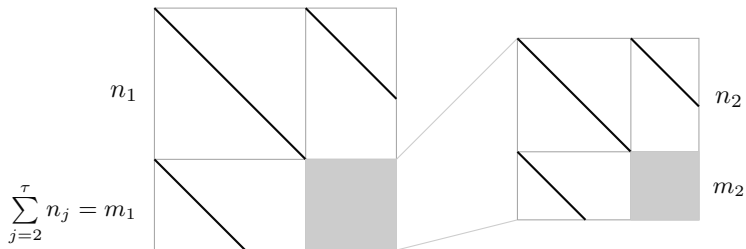
where $\{S^{(k)}\} \subset \mathbb{S}^n$ are symmetric and doubly stochastic with banded pattern



Doubly stochastic hierarchically banded (DSHB) factor

$$J = J_0 A_{\text{DSHB}} J_0,$$

where A_{DSHB} is symmetric, doubly stochastic, and hierarchically banded



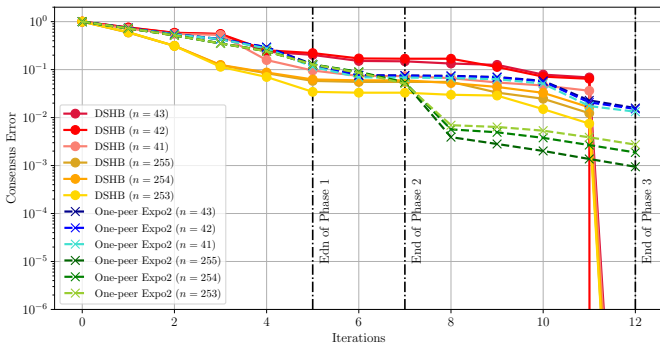
Numerical demonstration: decentralized average consensus

- decentralized average consensus iterations

$$x_i^{(k+1)} = W^{(k)} x_i^{(k)}, \quad \text{for } i = 1, \dots, n \text{ in parallel}$$

- we plot the consensus error

$$\Xi(k) = \frac{1}{n} \sum_{i=1}^n \|x_i^{(k)} - \bar{x}\|_2^2$$



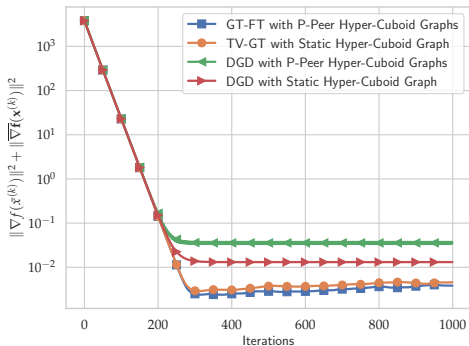
Numerical demonstration: decentralized optimization

consider the nonconvex optimization problem

$$\text{minimize} \quad \frac{1}{n} \sum_{i=1}^n \|A_i x - b\|^2 + \mu \sum_{j=1}^d \frac{x[j]^2}{1+x[j]^2}$$

- apply decentralized gradient descent (DGD) and gradient tracking (GT)
- static counterpart is built as the union of the graph sequence

$$E^{(\text{static})} = E^{(0)} \cup \dots \cup E^{(\tau-1)}$$



Summary

we study three classes of graph sequences with finite-time consensus

graph sequence	size n	# iterations
p -peer hyper-cuboids [NJYU'23]	any $n \in \mathbb{N}_{\geq 2}$	$n = \prod_{k=1}^{\tau} n_k$
SDS factor graphs [JNUY'24]	any $n \in \mathbb{N}_{\geq 2}$	$n = \sum_{k=1}^{\tau} n_k$
DSHB factor graphs [JNUY'24]	any $n \in \mathbb{N}_{\geq 2}$	$n = \sum_{k=1}^{\tau} n_k$

- finite-time consensus is achieved for **any** $n \in \mathbb{N}_{\geq 2}$
- takes into account **intra-cluster** and **inter-cluster** communications

Application to decentralized optimization

- reduced communication cost when used in existing decentralized methods
- **algorithm development:** more to expect ...